

- Regression
- Causation
- Correlation
- Discuss Quizzes/Answer Questions

6.0 Lesson Plan

- non-zero r does not imply a causal relationship.
 - r equals -1 if all points lie on a line with negative slope.
 - r equals 1 if all points lie on a line with positive slope.
 - r lies between -1 and 1, inclusive.
- Correlation coefficient r measures the strength of the linear association between X and Y values in a **scatterplot**.

6.1 More About Correlation

Some examples of scatterplots.

But sometimes, there might be a causal link. Hours of study is probably correlated with GPA, and it seems likely to be causal.

cause SAT, and Rob Zombie does not hurt GPA.
So it is hard to argue that correlation implies causation. GPA does not

affected by lifestyle.

- number of hours spent listening to Rob Zombie and GPA are both GPA and SAT scores are both affected by IQ.

Correlations are often high when some factor affects both X and Y .

6.2 Causation

What kinds of confounding factors did the tobacco lobby suggest?

The original link between smoking and lung cancer was an ecological correlation (Doll, 1955). The scatterplot showed the lung cancer rate against the proportion of smokers for 11 different countries.

Ecological correlations occur when X or Y or both is an average or a percentage for a group. Here causation is especially difficult to show.

What does a plot of the breast cancer rate against per capita fat consumption for various countries look like, and what conclusion should one draw?

Suppose one plotted GPA against SAT score for a set of individuals. Would this be an ecological correlation?

Is this an ecological correlation? What might be going on?
A study showed that cantons in Switzerland that had high literacy rates also had high suicide rates.

In multiple regression there are more than one kind of explanatory variable. For example, one might try to predict one's grade in statistics using both high school GPA and IQ.

independent variable, or the covariate.

- The explanatory variable is labelled X . This is sometimes called the dependent variable.
- The response variable is labelled Y . This is sometimes called the dependent variable.

Regression terminology:

6.3 Regression

Regression tries to fit the “best” straight line to the data. Specifically, it fits the line that minimizes the sum of the squared deviations from each point to the line, where deviation is measured in the vertical direction. Note: This does not measure deviation as the perpendicular distance from the point to the line.

- The mathematical model for regression assumes that:
- Each point (X_i, Y_i) in the scatterplot satisfies:
- $$Y_i = a + b X_i + \epsilon_i$$
- where the ϵ_i have a normal distribution with mean zero and (usually) unknown standard deviation.
- The errors ϵ_i have nothing to do with one another. A large error does not tend to be followed by another large error, for example.
 - The X_i values are measured without error. (Thus all the error occurs in the vertical direction, and we do not need to minimize perpendicular distance to the line.)

$$\frac{\sum_{i=1}^n X_i^2 - \bar{X}_2^2}{\sum_{i=1}^n X_i Y_i} = q \quad ; \quad \bar{X}q - \bar{Y} = \hat{a}$$

to zero, and solve. One finds that:

So take the derivative of $f(a, b)$ with respect to a and b , set these equal

$$\cdot \left[(\bar{X}q + \hat{a}) - \bar{Y} \right] \sum_{i=1}^n = f(a, b)$$

The sum of the squared vertical distances is

the squared vertical distance.

regression equation? We need to get the values that minimize the sum of the squared vertical distances in the regression equation?

How does one find the estimates \hat{a} and \hat{b} of the coefficients in the